

P-2213

## PATENT ABSTRACTS OF JAPAN

(11)Publication number : 11-272686

(43)Date of publication of application : 08.10.1999

(51)Int.Cl.

G06F 17/30  
G06F 17/27

(21)Application number : 10-070688

(71)Applicant : NIPPON TELEGR & TELEPH CORP  
<NTT>

(22)Date of filing : 19.03.1998

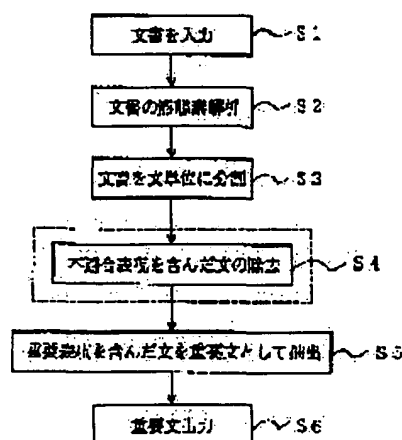
(72)Inventor : HORII MUNEYUKI  
MATSUOKA KOJI  
TAKAGI SHINICHIRO

(54) METHOD AND DEVICE FOR IMPORTANT DOCUMENT SENTENCE EXTRACTION AND RECORD MEDIUM WHERE IMPORTANT DOCUMENT SENTENCE EXTRACTING PROGRAM IS RECORDED

(57)Abstract:

PROBLEM TO BE SOLVED: To easily extract an important sentence from a document with high precision.

SOLUTION: An improper expression table wherein improper expressions are described and an important expression table wherein important expressions are described are prepared; and a morpheme analysis (S2) of an inputted document is taken, the analyzed document is divided (S3) into sentences, and sentences including improper expressions are removed (S4) from the document divided into the sentences by referring to the improper expression table. From the document from which the sentences including the improper expression have been removed, sentences including important sentences are extracted (S5) as important sentence by referring to the important expression table. Here, the process for removing the sentences including the improper expressions is omitted in some cases and the sentences including the important expression may be extracted as important sentence directly from the document divided into the sentences by referring to the important expression table.



## LEGAL STATUS

[Date of request for examination] 10.05.2001

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

(19)日本国特許庁 (JP)

(12) 公開特許公報 (A)

(11)特許出願公開番号

特開平11-272886

(43)公開日 平成11年(1999)10月8日

SI)Int.Cl.*		P I		識別記号	
G 0 6 F	17/30	G 0 6 F	15/401	3 2 0 A	
	17/27		15/20	5 5 0 E	
			15/38	D	
			15/40	3 7 0 A	

審査請求 未請求 請求項の数 5 O L (全 7 頁)

(21)出願番号	特開平10-70688	(71)出願人	00004228 日本電信電話株式会社 東京都千代田区大塚二丁目3番1号
(22)公開日	平成10年(1998)3月19日	(72)発明者	堀井 親之 東京都新宿区西新宿三丁目19番2号 日本 電信電話株式会社内
		(72)発明者	松岡 浩司 東京都新宿区西新宿三丁目19番2号 日本 電信電話株式会社内
		(72)発明者	高木 伸一郎 東京都新宿区西新宿三丁目19番2号 日本 電信電話株式会社内
		(74)代理人	弁護士 鈴木 誠

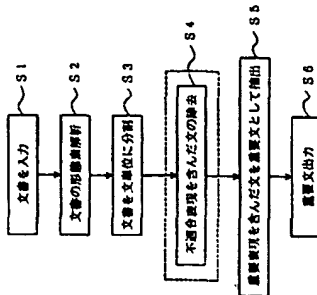
(54)【発明の名称】 文書重要文抽出方法、文書重要文抽出装置及び文書重要文抽出プログラムを記録した記録媒体

(57)【要約】

【課題】 文書中の重要文を、簡単に且つ高い精度で抽出することを可能にする。

【解決手段】 重要文として不適切な表現を記述した不適表現テールと、重要表現を記述した重要表現テールを用意し、入力された文書を形態素解析し (S2)、該形態素解析された文書を文単位に分割し (S3)、該文単位に分割された文書中から不適表現テールを参照し不適表現を含んだ文を取り除き (S4)、該不適表現を含んだ文を取り除いた後の文書中から重要表現テールを参照して重要表現を含んだ文を重要文として抽出する (S5)。ここで、場合によっては、不適表現を含んだ文を取り除く処理は省略し、文単位に分割された文書中から、重要表現テールを参照して直接重要表現を含んだ文を重要文として抽出することによい。

本発明の文書重要文抽出プログラムの処理フロー図



## 【特許請求の範囲】

【請求項1】 文書を入力し、該入力された文書を形態素解析し、該形態素解析された文書を文単位に分割し、該文単位に分割された文書中から重要表現を含んだ文を重要文として抽出することを特徴とする文書重要文抽出方法。

【請求項2】 請求項1記載の文書重要文抽出方法において、重要文として不適な表現を記述した不適表現テーブルを参照して、文書中から不適表現を含んだ文を取り除いた後の文書を重要文抽出の対象とすることを特徴とする文書重要文抽出方法。

【請求項3】 文書を入力する手段と、前記入力された文書を形態素解析する手段と、前記形態素解析された文書を文単位に分割する手段と、重要表現を記述した重要表現テーブルと、前記重要表現テーブルを参照して、前記文単位に分割された文書中から重要表現を含んだ文を重要文として抽出する手段を有することを特徴とする文書重要文抽出装置。

【請求項4】 文書を入力する手段と、前記入力された文書を形態素解析する手段と、前記形態素解析された文書を文単位に分割する手段と、重要文として不適な表現を記述した不適表現テーブルと、前記不適表現テーブルを参照して、前記文単位に分割された文書中から不適表現を含んだ文を取り除く手段と、重要表現を記述した重要表現テーブルと、前記重要表現テーブルを参照して、前記不適表現を含んだ文を取り除いた文書中から重要表現を含んだ文を重要文として抽出する手段を有することを特徴とする文書重要文抽出装置。

【請求項5】 文書から重要文を抽出するための文書重要文抽出プログラムを記録したコンピュータ読み取り可能な記録媒体であって、入力された文書を形態素解析する処理プロセスと、形態素解析された文書を文単位に分割する処理プロセスと、重要文として不適な表現を記述した不適表現テーブルを参照して、文単位に分割された文書中から不適表現を含んだ文を取り除く処理プロセスと、重要表現を記述した重要表現テーブルを参照して、不適表現を含んだ文を取り除いた文書中から重要表現を含んだ文を重要文として抽出する処理プロセスを有することを特徴とする記録媒体。

## 【発明の詳細な説明】

## 【0001】

【発明の属する技術分野】 本発明は、コンピュータによる文書処理に係り、詳しくは、文書の要約等を目的として文書の重要文を抽出する文書重要文抽出方法、文書重要文抽出装置、及びそのための文書重要文抽出プログラムを記録した記録媒体に関する。

## 【0002】

【従来の技術】 大量の情報、特にWWW (World Wide Web) の閲覧情報や電子メールといった電子的な情報

を効率良く閲覧し、必要な情報だけを取得するためには、文書の要約技術が必要不可欠である。

【0003】 従来、コンピュータによる文書処理において、文書を要約するために文書中から重要文を自動で抽出する技術としては、「文書中の単語頻度を利用する手法」、「文書構造を解析する手法」がよく用いられる。「文書中の単語頻度を利用する手法」は、大量の文書を統計処理することにより、文書中に現れる各単語を単語単位の出現頻度、文書単位の出現頻度を基にランキングし、そのランキングが高い単語を含む文を重要文として選択する手法である。一方、「文書構造を解析する手法」は、文書の構造から主題があるべきブロックを特定し、その中から重要文を選択したり、話題の導入、転換等を表す言語表現を指し示す重要文を選択する手法である。

## 【0004】

【発明が解決しようとする課題】 従来技術の「文書中の単語頻度を利用する手法」は、大量の文書データが用意できる場合には有効な方法となるが、反面、(1) 大量の文書データが必要である。(2) 必ずしもランキングが高い単語 (一般的に文書内出現頻度が高く、他文書における出現頻度が少ない単語) が重要文に含まれているとは限らない。(3) 該当単語の出現頻度が非常に高い場合、ほとんどの文に該当単語が含まれて多数の文が選択されてしまう、などの問題点がある。

【0005】 また、「文書構造を解析する手法」は、文書構造をうまく解析できれば効果的であるが、文書構造を正しく解析できなければならず、話し言葉のように解が困難な文書や、電子メールのように定型な構造を持たない文書を対象とする場合、重要文抽出までたどりつかない問題がある。

【0006】 本発明は、上記従来技術の問題点に鑑みてなされたもので、文書中から簡単かつ高い精度で、重要表現を含んだ文を重要文として抽出する文書重要文抽出方法、文書重要文抽出装置、及びそのための文書重要文抽出プログラムを記録した記録媒体を提供することを目的とする。

## 【0007】

【課題を解決するための手段】 上記目的を達成するため、本発明は、重要文として不適な表現と重要表現を記述したテーブルを用意し、文書中から不適表現を含んだ文を取り除いた後の文書を対象に、重要表現を含んだ文を重要文として抽出することを特徴とするものである。なお、場合によっては、不適表現を含んだ文を取り除く処理は省略してもよい。

## 【0008】

【発明の実施の形態】 図1は、本発明の文 重要文抽出方法の処理フロー図である。図1に示すように、本発明の方法では、文書を入力 (ステップS1)、該入力された文書を形態素解析し (ステップS2)、該形態素解析

された文書を文単位に分割し (ステップS3)、まず、重要文として不適な表現を記述した不適表現テーブルを参照して、文単位に分割された文書中から不適表現を含んだ文を取り除き (ステップS4)、次に、重要表現を記述した重要表現テーブルを参照して、不適表現を含んだ文を取り除いた文書中から重要表現を含んだ文を重要文として抽出し (ステップS5)、該抽出した重要文を出力する (ステップS6)。ここで、ステップS4の不適表現を含んだ文を取り除く処理は、場合によっては省略し、文単位に分割された文書中から、重要表現テーブルを参照して直接重要表現を含んだ文を重要文として抽出してもよい。

【0009】 図2は、本発明の文書重要文抽出装置の一実施形態のブロック図である。本文書重要文抽出装置は、文書入力装置10、形態素解析部21、単文分割部22、不適文除去部23及び重要文抽出部24の各処理機能を有する文書処理装置本体20、文書出力装置30、及び、単語辞書41、不適表現テーブル42、及び重要表現テーブル43を格納する記憶装置40などで構成される。この構成は、所謂コンピュータシステムで実現されるものである。ここで、文書入力装置10は、キーボード、イメージスキャナ、FD-やCD-ROM等のドライブ、あるいは、WWWや電子メールでの文書入力のための通信インタフェースの総称である。文書処理装置本体20は、実行するプログラムや所蔵データ等を格納するための内部メモリを有する所謂CPU本体である。記憶装置40はハードディスクなどである。

【0010】 図2において、文書入力装置10は、重要文抽出の対象となる文書の入力を行い、入力された文書を形態素解析部21に送出する。形態素解析部21は、文書入力装置10より受け取った文書を単語辞書41を参照して形態素解析を行い、この形態素解析された文書を単文分割部22に送出する。単文分割部22は、形態素解析部22より受け取った文書を文単位に分割する。ここでは従来と同様である。単文分割部22では、この文単位に分割された文書を不適文除去部40に送出する。

【0011】 不適文除去部23は、不適表現テーブル42を参照することにより、不適表現部22より受け取った文書中から不適表現を含んだ文を取り除く。不適表現テーブル42には、重要文として不適な表現が形態素情報、文位置情報として記述されている。不適文除去部23により、不適表現テーブル42に記述されたいずれかの表現を含んだ文全てが文書中から取り除かれる。不適表現を含んだ文を取り除いた文書を重要文抽出部24に送出される。

【0012】 重要文抽出部24は、重要表現テーブル43を参照することにより、不適文除去部23より受け取った文書中から重要表現を含んだ文を抽出する。重要表

現テーブル43には、重要表現が形態素情報、文位置情報として記述されている。重要文抽出部24により、重要表現テーブル43に記述されたいずれかの表現を含んだ文全てが重要文として文書中から抽出される。抽出された重要文は文書出力装置30に出力される。

【0013】 次に、具体例を用いて、本発明による文書から重要文を抽出する手法を説明する。

【0014】 文書入力装置10で、以下の電子メール文書が入力された場合を考える。

## &lt;入力&gt;

情須要情報通過 高本様

三田商事の竹石です。いつもお世話になっております。システム緊急会議の件でございます。

貴社に納入して頂いた在庫管理システムが今朝ダウンし、現在非常モードで稼働中です。主な状況は以下の通りです。

・バックアップデータがリストアできません。  
・年末時間と業務に支障があり、早急に対応を依頼しなければなりません。至急電話を頂けないでしょうか？本日15時から17時に三田商事 本社会議室で緊急対策会議を行いますので、ご出席をお願いいたします。申し訳ありませんが、よろしくお願いたします。

【0015】 形態素解析部21では、単語辞書41を照して上記入力文書を形態素解析する。たとえば、上記入力文書のうち、「三田商事の竹石です。」を形態素解析した結果は、図3のようになる。図3において、上段が原文及び分割点、中段が各単語の標準表記、下段が各単語の品詞を表す。文書入力装置10で入力された文書は全て形態素解析部21により、図3と同様に形態素解析される。

【0016】 単文分割部22では、形態素解析された入力文書を文単位に分割する。上記入力文 は、以下のように入文単位に分割される。ただし、以下では、説明の都合上、分割された文の順番に文書 を付与して示す。また、形態素情報は省略する。

1 情須要情報通過 高本様

2 三田商事の竹石です。

3 いつもお世話になっております。

4 システム緊急会議の件でございます。

5 貴社に納入して頂いた在庫管理システムが今朝ダウンし、現在非常モードで稼働中です。

6 主な状況は以下の通りです。

7 ・議室からアクセスしても初期画面が表示されません。

8 ・バックアップデータがリストアできません。

9 年末時間と業務に支障があり、早急に対応を依頼しなければなりません。

10 至急電話を頂けないでしょうか？

11 本日15時から17時に三田商事本社会議室で緊

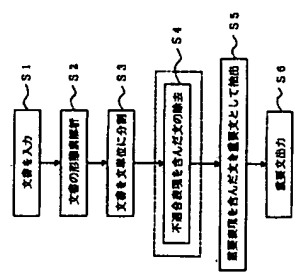


4 0 記憶装置  
4 1 単語辞

4 2 不適表現テーブル  
4 3 重要表現テーブル

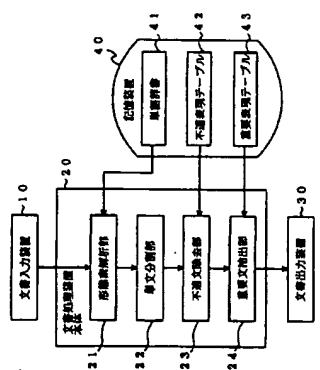
【図 1】

本発明の文書重要文抽出方法の処理フロー図



【図 2】

本発明の文書重要文抽出装置の一実施例の構成ブロック図



【図 3】

形態素解析結果の具体例

三田 / 同書 / の / 竹石 / です / 。 /  
三田 同書 の 竹石 だ  
題名欄(書名) 巻名欄 著者名欄(氏) 刊行所 句点

【図 4】

本発明に用いる不適表現テーブルの内容例

テーブル ID	文化圏情報	不適表現	不適表現	不適表現	品名
100	*	and	*	*	不適表現
101	*	not	*	*	不適表現
102	*	not	*	*	不適表現
103	*	not	*	*	不適表現

【図 5】

本発明に用いる重要表現テーブルの内容例

テーブル ID	文化圏情報	重要表現	重要表現	重要表現	品名
200	*	not	*	*	不適表現
201	*	not	*	*	不適表現
202	*	not	*	*	不適表現
203	*	not	*	*	不適表現
204	*	not	*	*	不適表現